

ASIST: Automatic Semantically Invariant Scene Transformation

Or Litany^{a,b}, Tal Remez^{a,b}, Daniel Freedman^a, Lior Shapira^c, Alex Bronstein^b, Ran Gal^c

^aMicrosoft Research, Haifa, Israel

^bElectrical Engineering Department, Tel Aviv University, Tel Aviv, Israel

^cMicrosoft Research, Redmond, WA, USA

Abstract

We present ASIST, a technique for transforming point clouds by replacing objects with their semantically equivalent counterparts. Transformations of this kind have applications in virtual reality, repair of fused scans, and robotics. ASIST is based on a unified formulation of semantic labeling and object replacement; both result from minimizing a single objective. We present numerical tools for the efficient solution of this optimization problem. The method is experimentally assessed on new datasets of both synthetic and real point clouds, and is additionally compared to two recent works on object replacement on data from the corresponding papers.

Keywords: Semantic segmentation, object recognition, random forest, Iterative Closest Point, Alternating minimization, Pose estimation, registration

1. Introduction

The problem we tackle in this paper is the transformation of 3D scenes. In particular, we are interested in the subclass of transformations which preserve *semantic invariance*: objects within the scene are to be replaced by other objects from the same class. Thus, a nightstand should be replaced by another nightstand, and not a packing box. While a particular packing box may be geometrically similar to the nightstand, it is semantically different and should therefore not be used as the replacement. Of course, to the extent possible, we would like to preserve geometric similarity as well; the nightstand should ideally be replaced by a nightstand with similar proportions, shape, position, and orientation. An example is given in Figure 1.

Semantically invariant scene transformation has a number of interesting applications. In the area of virtual reality, these transformations are useful for designing virtual scenes matching the underlying real scene in which the user is located. While not critical when the user is stationary, accurate object placement is crucial in any scenario in which the user moves throughout the scene. For example, semantic invariance means that in sitting on a virtual chair, the user is actually sitting on a real chair. Beyond preventing injury, this leads to a more realistic VR experience. In a different application, semantically invariant transformations may be used in the repair of point clouds acquired by the stitching together of many depth images, such as in [1]. These fused scans typically have occlusions, holes, and other artifacts, which could be mitigated by the replacement of scene objects with their pristine versions, based on CAD mod-

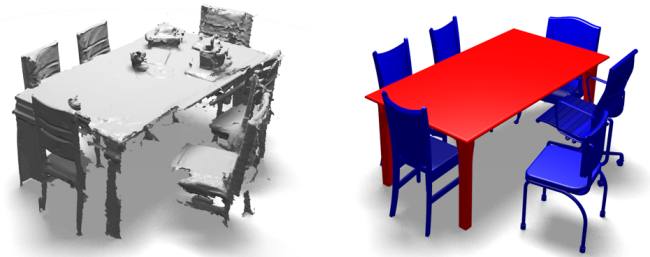


Figure 1: **Example output of the ASIST algorithm.** Left: point cloud of a scene acquired using the Kinect sensor. Right: ASIST output. Note that semantic invariance has been preserved: chairs replace chairs, and likewise for tables.

els. A third application involves mobile robotics systems. Collision avoidance – a necessary component of mobile robotics – is aided by more complete 3D data, and these systems would enjoy benefits similar to the point cloud repair scenario.

While semantically invariant scene transformation is not as standard a problem as, say, object detection or recognition, there have been a few studies in this direction in the last few years. Both Nan *et al.* [2] and Li *et al.* [3] present approaches to very similar problems, though the approaches themselves differ from ours. In Nan *et al.* [2], patches of the scene which are likely candidates for replacement are greedily added, and the resulting patch collection is matched against an object database. Li *et al.* [3] instead use collections of keypoints, which are then similarly matched against a database. By contrast, Gupta *et al.* [4] present a complex pipeline rather than a single algorithm. The pipeline involves many stages, including amongst others: contour detection, perceptual grouping, a convolutional neural network for object detection, instance segmentation, a second convolutional neural network for pose estimation, and registration. It is also noteworthy that [4] is aimed at RGB-D images,

Email addresses: t-orlita@microsoft.com (Or Litany),
t-talrem@microsoft.com (Tal Remez), danifree@microsoft.com
(Daniel Freedman), liors@microsoft.com (Lior Shapira),
bron@eng.tau.ac.il (Alex Bronstein), rgall@microsoft.com (Ran Gal)

rather than point clouds. Other works which are related, though perhaps not as closely, are SLAM++ of Salas *et al.* [5], and the Sliding Shapes approach of Song and Xiao [6].

Contributions. We refer to our approach to the problem as ASIST: Automatic Semantically Invariant Scene Transformation. ASIST is formulated as a single unified algorithm. Specifically, both the semantic labeling of the point cloud as well as the replacement of objects within the point cloud are presented as the solution to a single optimization problem. This contrasts most strongly with the work of Gupta *et al.* [4], which - while presenting very impressive results - is a complex system composed of many individual algorithms. We believe that the formulation of a single, unified algorithm for solving semantically invariant scene transformation is an important contribution in its own right.

As has already been noted, another difference from some of the prior work is in ASIST's focus on point clouds, as opposed to RGB-D images. The differences here are twofold. First, point clouds are missing the RGB component, which represents a potential advantage for several reasons. Low-power and cheap depth sensors usually do not possess an RGB sensor; moreover, even if they do, it is expensive to keep the RGB sensor on and operating all the time. Furthermore, the RGB quality is often poor, especially in the most relevant indoor low-light scenario. Finally, the mechanical limitations of such sensors make the accurate calibration of the RGB and depth cameras a challenging task. Second, point clouds have more geometric information than depth images. The acquisition of point clouds has become more common with the growth of systems such as Kinect Fusion [1] and Project Tango [7], so focusing on this data structure is a natural choice.

Paper Organization. The remainder of the paper is organized as follows. Section 2 presents related work in more detail. Section 3 presents the formulation of ASIST, focusing on the unified treatment of semantic labelling and object replacement. Section 4 evaluates ASIST on datasets of both synthetic and real point clouds, and comparison with the approaches of Nan *et al.* [2] and Li *et al.* [3]. Section 5 presents an overall discussion of the algorithm and results, while Section 6 concludes.

2. Related Work

In recent years several studies have explored problems beyond classical object detection or recognition. In Nan *et al.* [2] a dataset of deformable models is used to detect object instances in a scanned scene. A given scene is over-segmented into smooth patches, and a RANSAC-like algorithm is used to add patches with high classification scores. The greedy collection of patches is then compared to models in the database, which are used to remove outliers.

In Salas *et al.* [5] object recognition is integrated into a SLAM (simultaneous localization and mapping) pipeline. During the process of integrating new frames into the scanning volume, an existing method [8] of object detection is employed. The algorithm as presented is able to search for a relatively

small number of models, with limitations primarily due to real-time requirements and the strength of the GPU.

Song *et al.* [6] render a large dataset of models from multiple angles and train a linear SVM classifier for each, according to the exemplar SVM approach of [9]. Using a sliding window over a scene, they run each of the classifiers and select a small number of possible models, as well as a best matching pose for each model. Their method produces impressive results at the cost of slow performance. The sheer number of classifiers and windows to be tested is prohibitive.

Li *et al.* [3] preprocess a dataset to extract and cluster keypoints from 3D models. Keypoints are arranged into constellations, forming a shape descriptor which is used in real-time 3D scanning to insert model instances into the scene. We use this paper as well as that of Nan *et al.* [2], and the datasets therein, as references for comparison to our method.

Gupta *et al.* [4] use an existing [10] object detection and segmentation framework and focus on finding object matches in a dataset. They train a convolutional neural network (CNN) to provide a small number of pose and model hypotheses, manually picking five exemplars per category. They then use ICP to perform fine alignment, and a linear classifier to select from the different hypotheses for each object instance.

Mahabadi *et al.* [11] tackle a related problem: 3D reconstruction of scenes using learned semantic models. In particular, they introduce a framework which formulates the task of scene reconstruction as a volumetric multi-label segmentation problem. The key idea in their approach is to assign semantic class (including free-space) indicator variables for each voxel. Inspired by the classical crystallographic technique of Wulff Shapes, they describe how anisotropic surface regularization, which penalizes transitions between labels, can be derived from training data. This naturally continues their line of work from [12] and [13], in which they propose other alternatives to this binary label penalty term, such as learning preferred directions for inter-label transitions [12] or anisotropic regularization for larger (non-convex) object parts [13]. In contrast to our approach, a unary data term uses information about rays between the camera and each voxel; thus, their work is more applicable to the case of depth images, rather than fused 3D scans as required in our case. Furthermore, in [11] the semantic labels are assigned to convex model parts and not to whole objects; hence additional computation would be needed to apply this method for the purpose of semantic scene transformation.

Our algorithm uses a dataset of 3D models containing a variety of object classes. We make use of the LightField Descriptor of [14] to find subclasses within the dataset and select a number of exemplars from each category. This descriptor is based on the idea that two similar 3D models look similar from all viewing angles. Multiple orthogonal projections of each model are encoded using Zernike moments and Fourier descriptors. The similarity between two models is then defined as the minimal similarity over the various rotations which can be applied to align the models.

In recent years several datasets and benchmarks for 3D models and scene have emerged. ShapeNet [15] is a collection of CAD models containing 3 million models, of which 220,000

are labeled. ModelNet [16] consists of 600 categories, including 40 main household objects categories. In the accompanying paper [16], algorithms for object detection and 2.5D object completion are trained and tested on this dataset. The SunRGBD dataset and benchmark [17] provides a standard to evaluate and measure the success of scene understanding algorithms. The dataset is densely annotated in both 2D and 3D, using a combination of polygons and bounding boxes.

3. Automatic Semantically Invariant Scene Transformation

3.1. General Approach

Our goal is to take a point cloud representing a real world scene as input, and to transform it while remaining faithful to the semantics of the scene. More specifically, given a fixed set of object classes and a database of objects from these classes, we wish to:

1. recognize instances of objects within the scene belonging to these classes;
2. replace these instances with geometrically similar objects of the same class from the database, respecting original positions and orientations.

The ideal output of the algorithm would then be a semantically similar scene, but with new objects placed within the scene. The potential applications of such a technique have been highlighted in Section 1.

In order to arrive at a more unified treatment of the problem, which also yields superior results, we treat both problems *simultaneously*. That is, our goal is to solve both the semantic segmentation and replacement problems at the same time and in a consistent manner. In what follows, we show how to formulate the problem to this end.

3.2. Cell Classification

We are given a point cloud representing the scene. We voxelize the scene, and organize these voxels into larger structures called *cells*. A cell is defined to be a cubic collection of voxels; that is, it is a patch of $m \times m \times m$ voxels. In practice, we take $m = 9$. Note that much of the scene is empty, so it may appear that voxelization is a wasteful process; however, we only perform computations on occupied cells, i.e. cells containing points from the scene’s point cloud.

Each voxel is regarded as the center of a cell surrounding it. We train a random forest [18] which classifies cells according to one of n_c fixed classes, or “clutter”, to which class label 0 is assigned. The forest performs its classification using split functions which are a generalization of decision stumps and binary decisions. More specifically, we compute a *base feature* for each voxel within the cell; we use three types of base features, namely binary occupancy, the distance function of the voxel’s center from the point cloud, and the height of the cell as measured by the height of its center voxel (for further implementation details, refer to Section 3.6). Thus, the base features may be written as $\mathbf{h}_k \in \mathbb{R}^{M_k}$; in the case of the binary occupancy and distance function, $M_k = m^3$, whereas in the case of

the height $M_k = 1$. The forest’s split functions are then represented as triples (k, \mathbf{u}, τ) where $k \in \{1, \dots, K\}$, $\mathbf{u} \in \mathbb{R}^{M_k}$, and $\tau \in \mathbb{R}$; the decision is then based on the value of the binary variable $\mathbb{1}[\mathbf{u}^T \mathbf{h}_k > \tau]$, where $\mathbb{1}[\cdot]$ denotes the indicator function. A decision stump corresponds to \mathbf{u} with only a single non-zero value; a pairwise decision has two non-zero values; and so on. We describe the details of the feature choices we use in practice in Section 3.6.

This procedure yields forest scores for each voxel; we then assign forest scores to each point p in the point cloud \mathcal{P} via nearest neighbors (though any simple interpolation scheme will do). Thus, at the end of this process we have a collection of forest scores f_{cp} , where f_{cp} indicates the score of class c for point p . As is customary for random forests, the scores for a given point form a probability distribution, i.e. $f_{cp} \geq 0$ and $\sum_c f_{cp} = 1$.

3.3. Joint Semantic Segmentation and Object Replacement

We are given a dataset of objects which come from the n_c fixed classes. To reduce computational complexity, the objects from each class within the database are clustered into a set of groups, each of which is represented by an *exemplar*. We achieve this by first doing a rough clustering based on scale using k -medoids; within each new cluster, we then further sub-cluster based on the LightField Descriptor [14], again using k -medoids. At the end of the process, we have a collection \mathcal{E} of exemplars taken from all classes.

Our goal is now to decide which of these exemplars to insert into our scene, where to insert them, and in which pose. To this end, we define the variables w_{ep} , representing the weight of an exemplar e at point p . For a fixed point, these weights can be thought of as a probability distribution over which exemplar should be inserted at that point. Ideally, then, we would have $w_{ep} = 1$ for exactly one exemplar e , and is $w_{ep} = 0$ for the rest. In practice, we will simply require the probability distribution condition, that is

$$\sum_e w_{ep} = 1 \quad \text{and} \quad w_{ep} \geq 0$$

for every e and p . By convention, $e = 0$ corresponds to “clutter”, or non-object; that is, $w_{0p} = 1$ is an indication not to insert any object at point p . Note that this representation is similar to the one used in [19].

Properly choosing the weights w_{ep} effectively solves a soft version of the semantic segmentation problem: we can assign a soft/probabilistic label to each point in the scene according to which exemplar that point corresponds to. But in addition, we wish to solve the object replacement problem. Replacement entails deciding which exemplars to insert into the scene, and in which pose – comprising both position and orientation. Neither the identity of the exemplars to be included nor their pose follows directly from the soft weights w_{ep} .

To solve the object replacement problem, therefore, we add two new sets of variables. (1) The variables $v_e \in [0, 1]$ are “votes” for the exemplar e , and indicate whether an exemplar should be inserted. A positive vote $v_e > 0$ indicates the exemplar is to be inserted in the scene, while a vote of $v_e = 0$ implies

it should be left out. (2) The transformations T_e denote the pose in which a candidate exemplar should be inserted within the scene. We take T_e to belong to the set of rigid transformations (translations and rotations), though one could broaden this to including scaling or more exotic non-rigid transformations.

We propose to perform semantic segmentation and object replacement jointly, by minimizing an energy which is the sum of six terms, that is:

$$E(\{w_{ep}\}, \{T_e\}, \{v_e\}) = \sum_{i=1}^6 \lambda_i E_i(\{w_{ep}\}, \{T_e\}, \{v_e\})$$

The six individual energy terms are as follows:

$$\begin{aligned} E_1 &= \sum_{p,c} \left(f_{cp} - \sum_e A_{ce} w_{ep} \right)^2 \\ E_2 &= \sum_{p,e} w_{ep} D(x_p, T_e \mathcal{X}_e) \\ E_3 &= \sum_e \sum_{p_1, p_2} L_{p_1 p_2} w_{ep_1} w_{ep_2} \\ E_4 &= \sum_p \sum_e |w_{ep}|^\ell \\ E_5 &= - \sum_e \left(\sum_p w_{ep} \right) v_e \\ E_6 &= \sum_{e_1, e_2} Q(T_{e_1} \mathcal{X}_{e_1}, T_{e_2} \mathcal{X}_{e_2}) v_{e_1} v_{e_2} \end{aligned}$$

The minimization is performed subject to the following constraints:

$$\sum_e w_{ep} = 1 \quad \forall p \quad w_{ep} \geq 0 \quad \forall p, e \quad 0 \leq v_e \leq 1 \quad \forall e$$

Let us take each of these terms in turn.

- E_1 – *Semantic Data Term*. The constants A_{ce} are defined as follows: $A_{ce} = 1$ if exemplar e belong to class c and 0 otherwise, reflecting the assignment of exemplars to classes known a priori. Thus, the semantic data term tries to drive the weights to be faithful to the forest classifier; the sum of the weights over all of the exemplars in a particular class should match the output of the forest classifier for that class.
- E_2 – *Geometric Data Term*. x_p is the location of the point p , and \mathcal{X}_e is the point cloud which represents exemplar e . $T_e \mathcal{X}_e$ indicates the point cloud resulting from applying transformation T_e to each point in \mathcal{X}_e . Finally, D denotes an extrinsic distance between two sets of points. Thus, the geometric data term is a kind of weighted distance term. The goal is to match the scene points as best as possible to a given exemplar. Note that this term only pays attention to exemplars with non-zero weights. For the clutter exemplar $e = 0$, $D(x_p, T_e \mathcal{X}_e)$ is not well-defined; we replace it by a constant $D_{clutter}$ for all points p . The actual value of $D_{clutter}$ used is discussed in Section 4.

- E_3 – *Spatial Smoothness Term*. This term is a spatial regularization on the weights. In particular, we strive to smooth out the weights considered as a function of spatial location; that is, for each exemplar e , we wish to smooth out w_{ep} , considered as a function of p . We achieve this by using a Laplacian smoothing term; here \mathbf{L} is the Laplacian matrix over a weighted graph defined over the point cloud. More details are given in Section 3.6.
- E_4 – *Sparsity Term*. The parameter ℓ is chosen to be in the range $(0, 1)$. Thus, this term promotes sparsity on the weights. Recalling that for each pixel p the weights w_{ep} are normalized to sum to 1, this term encourages all of the mass to be placed on a single exemplar. In practice, we choose $\ell = 0.1$.
- E_5 – *Weight-Vote Agreement Term*. This term prefers to choose a large vote v_e for exemplar e when the sum of the weights for that exemplar, over all pixels, is large.
- E_6 – *Non-Collision Term*. $Q(\mathcal{X}_1, \mathcal{X}_2)$ denotes the binary overlap between two shapes \mathcal{X}_1 and \mathcal{X}_2 i.e. $Q(\mathcal{X}_1, \mathcal{X}_2) = 1$ if \mathcal{X}_1 and \mathcal{X}_2 overlap and 0 otherwise. Thus $Q(T_{e_1} \mathcal{X}_{e_1}, T_{e_2} \mathcal{X}_{e_2})$ measures the binary overlap between the exemplars e_1 and e_2 , in the poses T_{e_1} and T_{e_2} , respectively. Given that the votes v_e are constrained to be non-negative, this term encourages us to select only non-overlapping exemplars.

Having thus defined an energy which captures the idea of joint semantic segmentation and object replacement, we proceed to describe how to minimize the energy.

3.4. Simplifying the Energy

We begin by rewriting the energy more neatly in matrix-vector notation. Suppose that the number of points, exemplars, and classes are n_p , n_e , n_c , respectively. We form the vector $\mathbf{w} \in \mathbb{R}^{n_e n_p}$ by stacking the weights, as

$$\mathbf{w} = [w_{11}, \dots, w_{1n_p}, \dots, w_{n_e 1}, \dots, w_{n_e n_p}]^T.$$

We similarly stack the outputs of the random forest into a vector $\mathbf{f} \in \mathbb{R}^{n_e n_p}$ as

$$\mathbf{f} = [f_{11}, \dots, f_{1n_p}, \dots, f_{n_e 1}, \dots, f_{n_e n_p}]^T.$$

We denote $d_{ep}(T_e) = D(x_p, T_e \mathcal{X}_e)$, and form the vector of distances $\mathbf{d}(\mathbf{T}) \in \mathbb{R}^{n_e n_p}$ by

$$\mathbf{d}(\mathbf{T}) = [d_{11}(T_1), \dots, d_{1n_p}(T_1), \dots, d_{n_e 1}(T_{n_e}), \dots, d_{n_e n_p}(T_{n_e})]^T.$$

We keep the transformation (\mathbf{T}) dependence explicit as we will be optimizing over the transformations. Finally, we define n_p matrices $\mathbf{R}_p \in \mathbb{R}^{n_e \times n_e n_p}$, which pick out the entries of \mathbf{w} related only to point p . In other words,

$$\mathbf{R}_p \mathbf{w} = [w_{1p}, \dots, w_{n_e p}]^T$$

We define a similar set of sampling matrices for the forest entries \mathbf{f} , which we denote \mathbf{S}_p . Finally, we denote by $\mathbf{Q}(\mathbf{T})$ the $n_e \times n_e$ matrix with entries $Q(T_{e_1} \mathcal{X}_{e_1}, T_{e_2} \mathcal{X}_{e_2})$.

We may then rewrite the terms of the objective as

$$\begin{aligned} E_1 &= \mathbf{w}^T \left(\sum_p \mathbf{R}_p^T \mathbf{A}^T \mathbf{A} \mathbf{R}_p \right) \mathbf{w} - 2\mathbf{f}^T \left(\sum_p \mathbf{S}_p^T \mathbf{A} \mathbf{R}_p \right) \mathbf{w} \\ E_2 &= \mathbf{d}(\mathbf{T})^T \mathbf{w} \\ E_3 &= \mathbf{w}^T (\mathbf{I}_{n_e} \otimes \mathbf{L}) \mathbf{w} \\ E_4 &= -\|\mathbf{w}\|_\ell^\ell \\ E_5 &= -\mathbf{v}^T (\mathbf{I}_{n_e} \otimes \mathbf{1}_{n_p}^T) \mathbf{w} \\ E_6 &= \mathbf{v}^T \mathbf{Q}(\mathbf{T}) \mathbf{v} \end{aligned}$$

where \otimes is the Kronecker product; \mathbf{I}_k is the identity matrix of size $k \times k$; and $\mathbf{1}_k$ is the vector whose entries are all 1, of dimension k . Simplifying, we then have

$$E = \mathbf{w}^T \Psi_{ww} \mathbf{w} + \theta_w(\mathbf{T})^T \mathbf{w} + \xi_w \|\mathbf{w}\|_\ell^\ell + \mathbf{v}^T \Psi_{vv}(\mathbf{T}) \mathbf{v} + \mathbf{v}^T \Psi_{vw} \mathbf{w}$$

where

$$\begin{aligned} \Psi_{ww} &= \lambda_1 \sum_p \mathbf{R}_p^T \mathbf{A}^T \mathbf{A} \mathbf{R}_p + \lambda_3 \mathbf{I}_{n_e} \otimes \mathbf{L} \\ \theta_w(\mathbf{T}) &= -2\lambda_1 \left(\sum_p \mathbf{R}_p^T \mathbf{A}^T \mathbf{S}_p \right) \mathbf{f} + \lambda_2 \mathbf{d}(\mathbf{T}) \\ \xi_w &= -\lambda_4 \\ \Psi_{vv}(\mathbf{T}) &= \lambda_6 \mathbf{Q}(\mathbf{T}) \\ \Psi_{vw} &= -\lambda_5 \mathbf{I}_{n_e} \otimes \mathbf{1}_{n_p}^T \end{aligned}$$

The constraints may be simplified as follows:

$$\Gamma \mathbf{w} = \mathbf{1} \quad \mathbf{w} \geq \mathbf{0} \quad \mathbf{0} \leq \mathbf{v} \leq \mathbf{1}$$

where Γ is a $n_p \times n_e n_p$ matrix, whose p^{th} row is equal to $\mathbf{1}_{n_p}^T \mathbf{R}_p$.

3.5. Minimizing the Energy

Our problem is now

$$\min_{\mathbf{T}, \mathbf{w}, \mathbf{v}} E = \mathbf{w}^T \Psi_{ww} \mathbf{w} + \theta_w(\mathbf{T})^T \mathbf{w} + \xi_w \|\mathbf{w}\|_\ell^\ell + \mathbf{v}^T \Psi_{vv}(\mathbf{T}) \mathbf{v} + \mathbf{v}^T \Psi_{vw} \mathbf{w}$$

subject to

$$\Gamma \mathbf{w} = \mathbf{1} \quad \mathbf{w} \geq \mathbf{0} \quad \mathbf{0} \leq \mathbf{v} \leq \mathbf{1}$$

We use a technique based on alternating minimization. In particular, we cycle through subproblems with respect to \mathbf{T} , \mathbf{w} , and \mathbf{v} . We now detail how to solve these individual subproblems.

Minimization w.r.t. \mathbf{T} . Fixing \mathbf{w} and \mathbf{v} , the minimization w.r.t. \mathbf{T} reduces to

$$\min_{\mathbf{T}} \theta_w(\mathbf{T})^T \mathbf{w} + \mathbf{v}^T \Psi_{vv}(\mathbf{T}) \mathbf{v}$$

There are two conditions under which the second term, $\mathbf{v}^T \Psi_{vv}(\mathbf{T}) \mathbf{v}$ is equal to 0. The first condition is that the coefficient on the non-collision term $\lambda_6 = 0$. The second condition is that positive overlap between two exemplars implies that at most one is selected; that is, $Q(T_{e_1} \mathcal{X}_{e_1}, T_{e_2} \mathcal{X}_{e_2}) > 0 \Rightarrow v_{e_1} =$

0 or $v_{e_2} = 0$. As we shall see, in practice one of these two conditions generally holds. That is, as we describe in both Algorithm 1 and Section 3.6, for the initial iteration λ_6 is set to 0, so that the first condition is satisfied. In later iterations when $\lambda_6 > 0$, we observe empirically that the second condition generally holds, at least approximately. We use this as justification to ignore the second term $\mathbf{v}^T \Psi_{vv}(\mathbf{T}) \mathbf{v}$.

In this case, the optimization problem reduces to

$$\min_{\mathbf{T}} \theta_w(\mathbf{T})^T \mathbf{w}$$

The role of the transformations \mathbf{T} is somewhat obscured within the matrix-vector formulation. Going back to the original formulation, we can rewrite the above as

$$\min_{T_1, \dots, T_{n_e}} = \sum_{p,e} w_{ep} D(x_p, T_e \mathcal{X}_e)$$

It is easy to see that this minimization problem is separable, i.e. it decomposes into a separate minimization for each exemplar e . So for each e one needs to solve

$$\min_{T_e} \sum_p w_{ep} D(x_p, T_e \mathcal{X}_e) \quad (1)$$

Solving (1) for the optimal transformation T_e is the same as solving a weighted ICP problem. This can be done with standard techniques, e.g. [20].

Minimization w.r.t. \mathbf{w} . The main issue here is the sparsity term $\|\mathbf{w}\|_\ell^\ell$, which is not convex. However, since we are already performing an iterative optimization, it is natural to use the iterative reweighted least squares (IRLS) technique. This allows us to replace $\|\mathbf{w}\|_\ell^\ell$ with

$$\sum_e \eta_{ep} w_{ep}^2$$

where

$$\eta_{ep} = |w_{ep}^{(k-1)}|^{\ell-2} \quad (2)$$

and $w_{ep}^{(k-1)}$ are the optimal weights from the previous, i.e. $(k-1)^{th}$ iteration. (For example, see [21].) In this case, the energy becomes

$$\tilde{E} = \mathbf{w}^T \tilde{\Psi}_{ww} \mathbf{w} + \theta_w(\mathbf{T})^T \mathbf{w} + \mathbf{v}^T \Psi_{vv}(\mathbf{T}) \mathbf{v} + \mathbf{v}^T \Psi_{vw} \mathbf{w}$$

where

$$\tilde{\Psi}_{ww} = \Psi_{ww} + \text{diag}(\boldsymbol{\eta})$$

Now, fixing \mathbf{v} and \mathbf{T} , it is clear that the minimization of \tilde{E} w.r.t. \mathbf{w} is a convex quadratic program, i.e.

$$\min_{\mathbf{w}} \mathbf{w}^T \tilde{\Psi}_{ww} \mathbf{w} + (\theta_w(\mathbf{T}) + \Psi_{vw}^T \mathbf{v})^T \mathbf{w} \quad \text{s.t.} \quad \Gamma \mathbf{w} = \mathbf{1}, \quad \mathbf{w} \geq \mathbf{0} \quad (3)$$

Thus, one can solve this step for the global minimum w.r.t. \mathbf{w} (not the global minimum of the entire function, just of this step) using standard solvers.

Input: point cloud \mathcal{P} , set of exemplars \mathcal{E}

Output: set of exemplars $\mathcal{E}_{rep} \subset \mathcal{E}$ to insert into the scene, pose T_e for each exemplar $e \in \mathcal{E}_{rep}$

Initialization:

- set the sequence of coefficients $\{\lambda_6^{(i)}\}_{i=1}^{N_{out}}$
- evaluate Random Forest to get initial confidence per class for each point f_{cp}
- run Mean Shift to set initial exemplar positions
- set initial \mathbf{w} according to forest scores:
 $w_{ep} \leftarrow f_{cp} / (\# \text{ exemplars in class } c)$
- at every position run N_{ICP} weighted ICP's with different initial rotation around z -axis; keep result with smallest distance
- set initial $\mathbf{v} \leftarrow \mathbf{1}$

Iterative Minimization:

```

for  $i \leftarrow 1$  to  $N_{out}$  do
   $E^{(i)} \leftarrow$  energy function with coefficients  $\lambda_6 = \lambda_6^{(i)}$ 
  for  $j \leftarrow 1$  to  $N_{in}$  do
    registration step – compute  $\mathbf{T}$ : solve (1)
    for  $k \leftarrow 1$  to  $N_{IRLS}$  do
      set  $\boldsymbol{\eta}$  according to (2)
      segmentation step – compute  $\mathbf{w}$ : solve (3)
    end
  end
  voting step – compute  $\mathbf{v}$ : solve (4)
end
 $\mathcal{E}_{rep} \leftarrow \{e \in \mathcal{E} : v_e > 0, \sum_p w_{ep} \geq \text{threshold}\}$ 

```

Algorithm 1: ASIST algorithm.

Minimization w.r.t. \mathbf{v} . One can easily observe that in fixing \mathbf{w} and \mathbf{T} , the minimization w.r.t. \mathbf{v} is a quadratic program, that is

$$\min_{\mathbf{v}} \mathbf{v}^T \boldsymbol{\Psi}_{vv}(\mathbf{T}) \mathbf{v} + (\boldsymbol{\Psi}_{vw} \mathbf{w})^T \mathbf{v} \quad \text{s.t.} \quad \mathbf{0} \leq \mathbf{v} \leq \mathbf{1} \quad (4)$$

However, given that $\mathbf{Q}(\mathbf{T})$, the collision matrix, is not positive semidefinite, the QP will not, in general, be convex. Thus, while one may use standard solvers, this step will not in general yield the global minimum w.r.t. \mathbf{v} ; rather, a local minimum is all that can be guaranteed.

The Algorithm. The overall algorithm is summarized in Algorithm 1. The main structure of the iterations consists of three nested loops: the outer loop over i , in which the energy function's coefficients are adjusted; the middle loop over j , which contains the minimizations over \mathbf{T} (registration step), and \mathbf{v} (voting step); and an inner loop over k , which allows for the IRLS iterations necessary for minimization over \mathbf{w} (segmentation step).

There are several details in the initialization which have not yet been covered. We now proceed to discuss these and other implementation details.

3.6. Implementation Details

In this section we explain in details the algorithm pipeline. Please refer to Algorithm 1 for each of the relevant steps.

Sequence of Energy Coefficients. The energy function is, in general, non-convex in \mathbf{v} ; this is due to the indefinite nature of the collision matrix $\mathbf{Q}(\mathbf{T})$, and hence of the matrix $\boldsymbol{\Psi}_{vv}(\mathbf{T}) = \lambda_6 \mathbf{Q}(\mathbf{T})$.

To alleviate this non-convexity, we define a monotonically increasing sequence of coefficients $\{\lambda_6^{(i)}\}_{i=1}^{N_{out}}$ on the non-collision term E_6 , and corresponding sequence of energy functions $\{E^{(i)}\}_{i=1}^{N_{out}}$. Choosing $\lambda_6^{(1)} = 0$ we get an initial energy function which is convex in \mathbf{v} . Using a small non-collision coefficient λ_6 will generally result in weights w_{ep} which are “non-decisive”, in that we will not have $w_{ep} = 1$ for a single e ; rather, for a given point several exemplars will have positive weights. By gradually increasing the non-collision coefficient λ_6 , the decisiveness of the weights improves at the cost of increasing non-convexity.

We observe experimentally that using a moderate growth rate for the non-collision coefficients λ_6 increases the probability of converging to the correct results. This is further ameliorated by initializing the quadratic program with the solution for \mathbf{v} from the previous iteration.

Random Forest. Recall that in Section 3.2, we described the computation of K base features per cell, where each such base feature is a vector $\mathbf{h}_k \in \mathbb{R}^{M_k}$. There are two separate types of base features that are used:

1. *Scalar field features.* These include the occupancy and distance function features. In this case, the size of the feature $M_k = m^3$, where the cell is $m \times m \times m$.
2. *Scalar features.* We used only such feature, the height of the cell, as measured by the height of the center voxel of the cell. In this case, the size of the feature is $M_k = 1$.

In describing the split functions used, we will treat scalar field features as functions in three dimensions, i.e. $h_k(x, y, z)$ over all voxels (x, y, z) in the cell, with the understanding that we can convert from functional notation $h_k(\cdot)$ to the vector notation \mathbf{h}_k used in Section 3.2 simply by stacking. Thus, whereas the split functions were previously described as $\omega(k, \mathbf{u}, \tau) = \mathbb{1}[\mathbf{u}^T \mathbf{h}_k > \tau]$, for the sake of a simpler explanation, we will now describe those corresponding to scalar field features by

$$\omega(k, \mathbf{u}, \tau) = \mathbb{1} \left[\sum_{(x,y,z) \in \text{Cell}} u(x, y, z) h_k(x, y, z) > \tau \right].$$

In all experiments we obtained a prior on the label probability of each point in the point cloud using a random forest classifier trained to maximize the Shannon Entropy at each split. The split functions ω used in our experiments were designed so that some enable rotation invariant splits, while others offer more rotation selective ones. In what follows, the z -direction is the vertical direction, measuring height off the ground.

- *Height:* Let h_3 be the height of the center voxel of the cell. Then $\omega = \mathbb{1}[h_3 > \tau]$.

- *Rotation Invariant*: Uses a dot product between a radially symmetric weight vector and the cells feature vector. $\omega = \mathbb{1}[\sum_{(x,y,z) \in \text{Cell}} u(x,y,z)h_k(x,y,z) > \tau]$ where u is a randomly generated function with values that are rotationally symmetric for rotations around the z -axis.
- *Box*: Sums the occupancy or the distance function in a randomly generated box. Thus, $\omega = \mathbb{1}[\sum_{(x,y,z) \in B} h_k(x,y,z) > \tau]$ where B is a box with a random location and size within the relevant cell, and $h_k(x,y,z)$ is either the occupancy or distance function value of the voxel at location (x,y,z) .
- *Horizontal Slab*: Sums the occupancy or distance function over all 9×9 voxels at a single z -value. The slab can be written as $R(z_0) = \{(x,y,z) \in \text{Cell} : z = z_0\}$. Then the split function is $\omega = \mathbb{1}[\sum_{(x,y,z) \in R(z_0)} h_k(x,y,z) > \tau]$, where $z_0 \in \{-4, -3, \dots, 4\}$ is selected randomly.
- *Pixelwise Values*: Linear combination of the occupancy or the distance function values of up to three voxels $\omega = \mathbb{1}[\sum_{i=1}^3 a_i h_k(x(v_i), y(v_i), z(v_i)) > \tau]$ where the voxels v_i are selected randomly within the cell limit and $a_i \in \{-1, 0, 1\}$ is also selected at random.

For all split functions the value of τ is selected randomly within the feasible bounds of the training data reaching the node.

Mean Shift for Initial Exemplar Positions. Initial exemplar locations are obtained by running a fast version of a weighted Mean Shift [22] algorithm on the point cloud after it has been projected onto the xy -plane. The algorithm is run once for each object class; during the run for object class c , the weight assigned to point p is taken to be f_{cp} , the random forest output for that class. The modes found by Mean Shift therefore tend to be at locations where the forest assigned high probabilities to the class in question.

The bandwidth of the Mean Shift algorithm is set differently for each class, based on a rough estimate of the object size in that class. As a post-processing step, we merge modes that are closer than 1.5 times the bandwidth. The modes returned for a given class are then assigned as the starting positions for each exemplar representing that class.

The next step in the initialization process is to determine the starting orientation of each exemplar. This is done by initializing each exemplar in 8 equally spaced rotations around the z -axis, running an ICP algorithm to finely tune the position and orientation of the exemplar, and then choosing the best fitted exemplar out of the 8. In cases in which the best exemplar is farther than a predetermined threshold all 8 candidates are removed from consideration.

Laplacian Operator \mathbf{L} for Spatial Smoothness Term E_3 . Given the point cloud \mathcal{P} we construct an undirected weighted graph $G(\mathcal{P}, \mathcal{F}, \mathbf{\Omega})$, where the edge set \mathcal{F} is constructed using k -nearest neighbors. The edge weights are given by $\Omega_{pq} = e^{-\|x_p - x_q\|^2 / 2\sigma^2}$ for $(p, q) \in \mathcal{F}$ and are zero otherwise. We use the random walk version of the normalized graph Laplacian of G [23] where $\mathbf{L} =$

$\mathbf{I} - \mathbf{\Delta}^{-1}\mathbf{\Omega}$, and $\mathbf{\Delta} = \text{diag}(\sum_q \Omega_{pq})$. In all of our experiments we used $k = 10$ nearest neighbors and $\sigma = 5$.

Output Exemplar Filtering. Recall from Section 3.3 that our proposed criterion for inserting exemplars was based entirely on their vote: if $v_e > 0$ exemplar e would be inserted, and otherwise it would not. In practice, this works quite well most of the time, as the algorithm generally produces a set of non-overlapping exemplars.

However, it may sometimes be the case that an exemplar is assigned in error to a small number of isolated points. This rarely happens since such points are usually classified as clutter by the forest, and thus Mean Shift finds no corresponding mode. But in the unusual case that a mode is found, then even if the vote v_e corresponding to that mode is small, it is still positive. Thus, the original criterion requires that we retain this exemplar, which is problematic.

To handle such cases we apply a straightforward filter as a final step, where we keep only exemplars e for which the vote v_e is positive, and at the same time the aggregation over the weights of the exemplar $\sum_p w_{ep}$ exceeds a small predetermined threshold. We set the threshold to be 0.1 in all our experiments.

3.7. Speeding Up the Algorithm

There are three steps to our algorithm, corresponding to the three sets of variables we minimize with respect to. Each of these steps must be repeated several times, as shown in Algorithm 1.

Minimizing with respect to the transformations T_e does not pose a speed problem, as weighted ICP may be solved quickly (see for example [24]). The issue is the two quadratic programs we must solve, one with respect to \mathbf{w} and the other with respect to \mathbf{v} . However, these two QPs are quite different in scale. The QP with respect to \mathbf{w} is of size $n_e n_p$, while the QP with respect to \mathbf{v} is of size n_e . Rough orders of magnitude for these two variables are $10^5 - 10^6$ for n_p and 10^2 for n_e . Thus, the QP for \mathbf{v} can be solved with a standard solver quite quickly; while the QP for \mathbf{w} , which will be of size $n_e n_p \approx 10^7 - 10^8$, is considerably slower to solve in a standard fashion. We thus propose two independent techniques to speed up this QP.

Subspace Parameterization. Let us denote by $\mathbf{w}_e \in \mathbb{R}^{n_p}$ the vector $[w_{e,1}, \dots, w_{e,n_p}]^T$. This is the vector of weights corresponding to just the exemplar e , over all points. We may therefore think of \mathbf{w}_e as a scalar field or function on the volume.

The idea is to represent this scalar function in a more parsimonious fashion. Thus, we use an expansion in terms of basis functions. A natural set of basis functions is provided by the Laplacian operator \mathbf{L} , which we already use in the smoothing term. We take the basis functions to be the eigenfunctions of \mathbf{L} , corresponding to the smallest n_b eigenvalues. (Recall that the smaller the eigenvalue of the Laplacian, the smoother the function.) We denote these functions by $\{\phi_i\}_{i=1}^{n_b}$, where each $\phi_i \in \mathbb{R}^{n_p}$, and the collection is given by the matrix whose columns are the individual functions, i.e. $\mathbf{\Phi} \in \mathbb{R}^{n_p \times n_b}$. In practice, we find $n_b = 30$ Laplacian basis functions suffice.

In this case, we can represent the function \mathbf{w}_e as $\Phi\alpha_e$, for $\alpha_e \in \mathbb{R}^{n_b}$. Stacking the coefficients of each exemplar to get a vector $\alpha \in \mathbb{R}^{n_b n_e}$, we have that

$$\mathbf{w} = (\mathbf{I}_{n_e} \otimes \Phi)\alpha \equiv \hat{\Phi}\alpha$$

And thus, we may represent the entire vector \mathbf{w} with only $n_b n_e$ values. Given that $n_b = 30$ in practice, this is a huge complexity reduction of 3-4 orders of magnitude.

The energy we need to minimize thus becomes

$$\begin{aligned} \tilde{E} &= \mathbf{w}^T \tilde{\Psi}_{ww} \mathbf{w} + (\theta_w(\mathbf{T}) - \Psi_{vv}(\mathbf{T})\mathbf{v})^T \mathbf{w} \\ &= \alpha^T (\hat{\Phi}^T \tilde{\Psi}_{ww} \hat{\Phi}) \alpha + [(\theta_w(\mathbf{T}) - \Psi_{vv}(\mathbf{T})\mathbf{v})^T \hat{\Phi}] \alpha \\ &\equiv \alpha^T \Psi_{\alpha\alpha} \alpha + \theta_\alpha(\mathbf{T})^T \alpha \end{aligned}$$

The two set of constraints are now represented in terms of the coefficients α as

$$\Gamma \hat{\Phi} \alpha = \mathbf{1} \quad \hat{\Phi} \alpha \geq \mathbf{0}$$

It is not clear if the equality constraints are even still feasible, given that we are dealing with a much smaller number of variables. It turns out they are still feasible, which we now show. To do so, it will be easier to go back to the original non-vectorized formulation. These constraints were $\sum_e w_{ep} = 1$ for all p . But this means that

$$\sum_e \mathbf{w}_e = \mathbf{1}_{n_p} \Rightarrow \sum_e \Phi \alpha_e = \mathbf{1}_{n_p} \Rightarrow \Phi \left(\sum_e \alpha_e \right) = \mathbf{1}_{n_p}$$

Thus, the question becomes: does $\Phi \mathbf{z} = \mathbf{1}$ have a solution? The answer is yes, due to the special properties of the Laplacian matrix, and its eigenfunctions. Let the eigenvalues of the Laplacian be denoted $\{\mu_i\}$. Then a result from spectral graph theory [23] shows that if β is given by

$$\beta_i = \begin{cases} \|\phi_i\|_1 & \text{if } \mu_i = 0 \\ 0 & \text{if } \mu_i > 0 \end{cases}$$

then $\Phi \beta = \mathbf{1}$, and it is the unique solution. This means that the equality constraints can be converted to

$$\sum_e \alpha_e = \beta \quad \Rightarrow \quad (\mathbf{1}_{n_e} \otimes \mathbf{I}_{n_p}) \alpha = \beta$$

which quite clearly have a multiplicity of solutions.

Thus, the QP becomes

$$\min_{\alpha} \alpha^T \Psi_{\alpha\alpha} \alpha + \theta_\alpha(\mathbf{T})^T \alpha$$

subject to

$$(\mathbf{1}_{n_e} \otimes \mathbf{I}_{n_p}) \alpha = \beta \quad \hat{\Phi} \alpha \geq \mathbf{0}$$

At first, this looks very reasonable: the number of variables is $n_b n_e$, which is quite manageable; and the equality constraints have been taken care of. The problem is the inequality constraints, specifically, the number of such constraints. Recall that these constraints derive from the non-negativity constraint on each weight, i.e. $w_{ep} \geq 0$; thus, there are still $n_e n_p$ such constraints. Since the complexity of quadratic programming depends on both the number of variables and the number of constraints, we will still have a slow algorithm.

Squared Weights. As a result, we must reduce the number of inequality constraints. We ask the following question: what would happen if we dropped these non-negativity constraints? Examining the various terms of the energy, we see that neither the semantic data term E_1 , nor the spatial smoothness term E_3 , nor the sparsity term E_4 would have any particular incentive to choose negative weights. Furthermore, the weight-vote agreement term E_5 would suffer significantly from choosing negative weights; and the non-collision term E_6 has no dependence on the weights. Thus, the only problematic term is the geometric data term $E_2 = \mathbf{d}(\mathbf{T})^T \mathbf{w}$. Quite clearly, this term would be minimized by choosing w as negative as possible.

To deal with this problem, we adopt the following simple workaround. We change the geometric data term E_2 , to be the following:

$$E'_2 = \sum_{p,e} w_{ep}^2 D(x_p, T_e \mathcal{X}_e) = \mathbf{w}^T \text{diag}(\mathbf{d}(\mathbf{T})) \mathbf{w}$$

That is, we replace the weights in the weighted ICP with the *squared weights*. Now, choosing large negative weights will pose a disadvantage. In fact, even small negative terms would tend to raise E'_2 , when coupled with the constraint $\sum_e w_{ep} = 1$; this is because if one term is negative, that means that the others must sum to more than 1.

Making this change yields the following changes in the problem parameters:

$$\Psi'_{ww}(\mathbf{T}) = \tilde{\Psi}'_{ww} + \lambda_2 \text{diag}(\mathbf{d}(\mathbf{T})) \quad \theta'_w(\mathbf{T}) = \theta_w(\mathbf{T}) - \lambda_2 \mathbf{d}(\mathbf{T})$$

and the inequality constraints are removed. That is, taking $\Psi'_{\alpha\alpha}(\mathbf{T}) = \hat{\Phi}^T \tilde{\Psi}'_{ww}(\mathbf{T}) \hat{\Phi}$ and $\theta'_\alpha(\mathbf{T}) = \theta'_w(\mathbf{T}) - \Psi_{vv}(\mathbf{T})\mathbf{v}$ we solve

$$\min_{\alpha} \alpha^T \Psi'_{\alpha\alpha}(\mathbf{T}) \alpha + \theta'_\alpha(\mathbf{T})^T \alpha \quad \text{s.t.} \quad (\mathbf{1}_{n_e} \otimes \mathbf{I}_{n_p}) \alpha = \beta \quad (5)$$

In this case, the QP can actually be solved as a linear system, as it is the result of minimizing a positive semi-definite quadratic form subject to a linear system. This yields a very fast algorithm in practice.

The Accelerated Algorithm. The faster algorithm is identical to Algorithm 1, with one change: the segmentation step now involves solving (5) for α instead of (3) for \mathbf{w} .

4. Results

To evaluate the performance of ASIST, we conducted several experiments using a variety of datasets including both our own synthetic and scanned scenes as well as scenes acquired by Nan *et al.* [2] and Li *et al.* [3]. Overall our datasets contain scans acquired by four different types of sensors: Kinect v1, Kinect v2, Mantis Vision, and Google Tango.

This section is organized as follows. In Section 4.1 we detail the parameters of ASIST, such as the number of iterations, weights for the different energy terms, etc. In Section 4.2 we evaluate ASIST on a dataset of scenes containing synthetic models obtained from ModelNet [17]. We evaluate our performance both quantitatively and qualitatively, and discuss success

and failure cases. In 4.3 we introduce a small scanned dataset containing scenes we acquired using Kinect v2 and Tango, and discuss the performance of ASIST on these scenes. Sections 4.4 and 4.5 include comparison with two different algorithms from the recent literature [3, 2]. Since there was no available implementation for either method, we ran ASIST on scanned datasets supplied by the authors on which they reported their performance. In both cases we found ASIST to give comparable results.

4.1. Experimental Settings

General Parameters. In all experiments we ran our algorithms for $N_{out} = 5$, $N_{in} = 2$, and $N_{IRLS} = 5$ iterations. The $\|\cdot\|_\ell$ sparsity inducing norm was chosen as $\ell = 0.1$. The energy term coefficients were taken to be $\lambda_3 = 100$, $\lambda_4 = 10$, and $\lambda_5 = 1$. As was described in Section 3.6, λ_6 was increased at each outer iteration, taking on values of 1, 5, 10, 10^2 , 10^3 , 10^9 . The rest of the parameters varied by experiment, and were set as shown in Table 1.

Forest Training. In all experiments we trained a forest with 9 trees of depth 10. At each node a pool of a 1,000 random split functions were generated and the one maximizing the Shannon information gain was selected. Training was stopped if the information gain was below 0.05 or if less than 30 samples reached the node. Each tree was trained of a random set of 6×10^5 cells selected out of models from ModelNet [17].

The average execution time per scene using unoptimized MATLAB code on an Intel Xeon E5620 2.4GHz is about 10 minutes.

	λ_1	λ_2	$D_{clutter}$	Voxel[cm]
Synthetic scenes	1	1	10	7.5
Data of Nan <i>et al.</i> [2]	10	1	20	2.5
Data of Li <i>et al.</i> [3]	10	10	10	2.5

Table 1: **Parameter settings.** Reported are only the parameters that vary across experiments. See discussion in the text.

4.2. Synthetic Dataset

Description of the Dataset. As an initial experiment, we evaluated the performance of ASIST on a synthetic benchmark. The benchmark comprised 30 scenes each containing a random collection of objects taken from the test portion of the ModelNet [17] dataset (recall that ASIST was trained on the non-overlapping training set of ModelNet). The objects were positioned at random non-overlapping locations within the scene. Objects from the following five classes were used in our benchmark: chair, table, toilet, sofa, and bed. Each scene contained either one or two objects from each class. The scene thus consists of a mesh, and is annotated additionally with a set of ground truth bounding boxes and their class labels. Note that the bounding boxes are not axis-aligned; rather, they are aligned tightly to the objects they surround.

Evaluation criteria. In all experiments we computed precision and recall using the bounding boxes; an overlap is considered to have occurred if the Intersection over Union (IoU) score, defined as the ratio between the intersection and union volumes of the bounding boxes [25], is greater than 0.25. We define two measures of relevance and, consequently, two types of precision and recall: *semantic*, for which the replaced and the replacing objects are deemed relevant if they have matching class labels, and *geometric*, which simply looks for overlap disregarding the class labels of the objects. Semantic precision and recall is reported on a per-class basis, whereas geometric precision and recall is given as an aggregate for the entire dataset.

We intend to release our synthetic benchmark to the research community.

Results. A quantitative evaluation of the performance of ASIST on the synthetic benchmark is given in Table 2. For each of the five classes, the semantic precision, recall, and F_1 scores are reported. Note that the semantic precisions vary between 0.91 and 1, with a mean precision of 0.97 across all classes; while the semantic recalls vary between 0.94 and 0.98, with a mean recall of 0.96 across all classes. The corresponding $F_1 = 2 \cdot \frac{\text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}}$ scores range between 0.94 and 0.99, with a mean value of 0.96. In addition, the geometric precision and recall (which do not account for class labels) are 1 and 0.99, with the geometric F_1 score of 0.99.

In our experiments, ASIST is executed for $N_{out} = 5$ outer loop iterations. We compare our results with those for which ASIST is run for $N_{out} = 1$. In this setting we assign a very high value for λ_6 . Given that the number of inner iterations is rather small ($N_{in} = 2$), $N_{out} = 1$ entails relatively little interaction between the segmentation, registration, and voting steps. Thus, comparing these two different settings of ASIST gives an indication of the importance of the unified approach to semantic segmentation and object replacement on which the ASIST algorithm is based. Several example scenes are shown in Figure 4.2. As can be observed from these examples, ASIST succeeds in properly identifying both the class and pose of the scene objects, and in finding a geometrically close substitute.

It is also evident from Figure 4.2 that the number of iteration affects the performance of ASIST. It can be seen in first example presented in Figure 4.2 that while ASIST with $N_{out} = 5$ returns a perfect result, ASIST with $N_{out} = 1$ replaces a table with a toilet (bottom right part of the scene), completely misses a chair (center of the scene), and places the bed at a wrong orientation (top right corner of the scene). In the second example ASIST with $N_{out} = 1$ replaces a toilet with a chair (center of the scene), completely misses a chair (top left corner of the scene), and once again positions the bed at a wrong orientation (top right corner of the scene). In the third example it replaces a sofa with a bed (bottom part of the scene).

Failure cases of the algorithm are shown in Figure 4.2. The example in the top row shows two different kinds of errors. First, a sofa in the scene is replaced by two chairs (top left of the scene). Second, both beds (in yellow, top right of the scene) are roughly recovered, but their pose is incorrect. In the example in the bottom row, a bed is replaced with a table (in

Measure \ Class	bed	chair	sofa	table	toilet	geo
Prec. ($N_{out} = 5$)	0.98	0.96	1	0.91	1	1
Prec. ($N_{out} = 1$)	0.97	0.93	0.98	0.91	0.97	0.99
Rec ($N_{out} = 5$)	0.95	0.94	0.96	0.98	0.98	0.99
Rec ($N_{out} = 1$)	0.83	0.88	0.91	0.78	0.88	0.89
F_1 ($N_{out} = 5$)	0.96	0.95	0.98	0.94	0.99	0.99
F_1 ($N_{out} = 1$)	0.9	0.91	0.95	0.84	0.92	0.94

Table 2: **Synthetic Benchmark Results.** Semantic and geometric precision, recall and F_1 scores for each of the 5 classes over the 30 benchmark scenes.

red, bottom right); the sofa at the bottom left is replaced with a considerably smaller sofa; and a rectangular table (top middle) is replaced with circular table. These sorts of errors are characteristic of ASIST’s failure modes.

The precision, recall, and F_1 scores for $N_{out} = 1$ are reported in Table 2. The mean semantic precision across classes for a single iteration is 0.95; when compared to 0.97 for five iterations, this demonstrates a modest improvement for the full algorithm. The improvement is much clearer in examining the recall scores, where the mean semantic recall across for a single iteration is 0.85, compared to 0.96 for five iterations.

Thus, we may conclude that the unified approach of ASIST to semantic segmentation and object replacement is indeed important; as it provides crucial improvements to the algorithm’s performance.

4.3. Fused Scans

Description of the Dataset. As explained in previous sections, the ASIST algorithm is designed to operate on point clouds obtained by a fusion of scans from multiple directions. To evaluate its performance in these scenarios, we collected several such fused scans using two types of sensors: Kinect v2 based on time-of-flight technology, and Tango based on triangulation [7]. We built a small collection of six indoor scenes, four of which were obtained using Kinect Fusion [1] and two using Tango. The scenes are presented in the leftmost columns of Figures 4.3-4.3 where it can be seen that even though the acquisition was done from multiple viewpoints, it consists of many holes, mainly due to occlusions. Data acquired using Tango is also much noisier, with approximately half the average point density of the scenes acquired with Kinect v2. The dataset consists a total of 29 chairs, 9 tables and 2 sofas.

Pipeline Illustration. In order to get a better understanding of the different steps of ASIST, we present intermediate results on the “dining table” scene shown in the top row of Figure 4.3. Figures 4.3 and 4.3 show intermediate outputs of the algorithm pipeline. The semantic segmentation is presented in Figure 4.3. Each point in the scene’s point cloud is assigned a semantic class confidence, which is just the sum of the weights of all exemplars belonging to the same class. The points are colored according to a heat-map, where red means high confidence and blue signifies low confidence. The leftmost column depicts the labeling achieved by the random forest at initialization for the chair class (top row) and for the table class (bottom row). In this example we used the six-class forest (clutter, chair, toilet,

bed, sofa and table) but since the scene contains only chairs and a table, for the sake of clarity we present only the relevant weights. It can be seen that the classifier succeeds in locating unique characteristics of the class such as the height of the plane in the table case and legs and backrest of the chairs. However, this results is far from being accurate enough for performing the scene transformation.

Figure 4.3 column (b) shows the result of using the spectral basis representation. Other than having the obvious benefit of reducing the number of optimization variables significantly, this low-pass filtering type of action has the effect of smoothing out the per-point labeling, thereby assisting the Mean Shift algorithm to avoid getting stuck in local modes. Finally, column (c) shows the final labeling achieved by ASIST. It can be seen than the objects points assume a much higher value of their corresponding class confidence. This demonstrates the power of the alternations done by the algorithm, i.e. segmentation contributes to better registration and vice versa. While this intermediate result is clearly not prefect, it is sufficiently accurate in order to perform a perfect replacement.

The evolution of exemplars though the course of $N_{out} = 5$ iterations is shown in Figure 4.3. The scene is shown in transparent gray from a top viewpoint and each exemplar is represented as a full circle located at its bounding box’s center. Each circle is assigned a color according to its class membership, using the same color code as described in 4.2. Column (a) depicts all the exemplar locations at initialization, i.e. after each exemplar has been placed at its corresponding class’ Mean Shift mode, and was then registered to the scene using several weighted ICP’s, each with a different initial rotation around z -axis. The location shown is of the exemplar position as output from the ICP that resulted in the smallest one-sided Hausdorff distance between the exemplar and the scene. Columns (b)-(d) show the locations of the exemplars e for which the vote v_e is positive and the aggregation of weights exceed a small threshold (this is the same criterion denoted by ϵ_{rep} in Algorithm 1).

Figure 4.3 demonstrates nicely how the large number of initial candidates is reduced at each step until we are left with only a final subset of exemplars which are both non-overlapping and have high confidence values. Increasing the value of λ_6 in a moderate fashion results in gradually removing the exemplars with low confidence while keeping the more ambiguous ones even if they are overlapping, thus leaving the hard selection for later steps where the best fitting exemplar presents a more significant advantage as compared to the others. Observe how after a single iteration the candidate exemplars belonging to the bed class have already been excluded, and after the third iteration only exemplars from the correct classes remain for ASIST to choose from. This gives good intuition for why ASIST tends to converge in practice to the correct solution even though it is not a convex problem; it is because the voting step (see Algorithm 1) is initialized close to the global solution and its non-convex term, $\lambda_6 E_6$ (the non-collision term), increases gradually.

Finally, the transformed scene is shown in the top row of Figure 4.3.

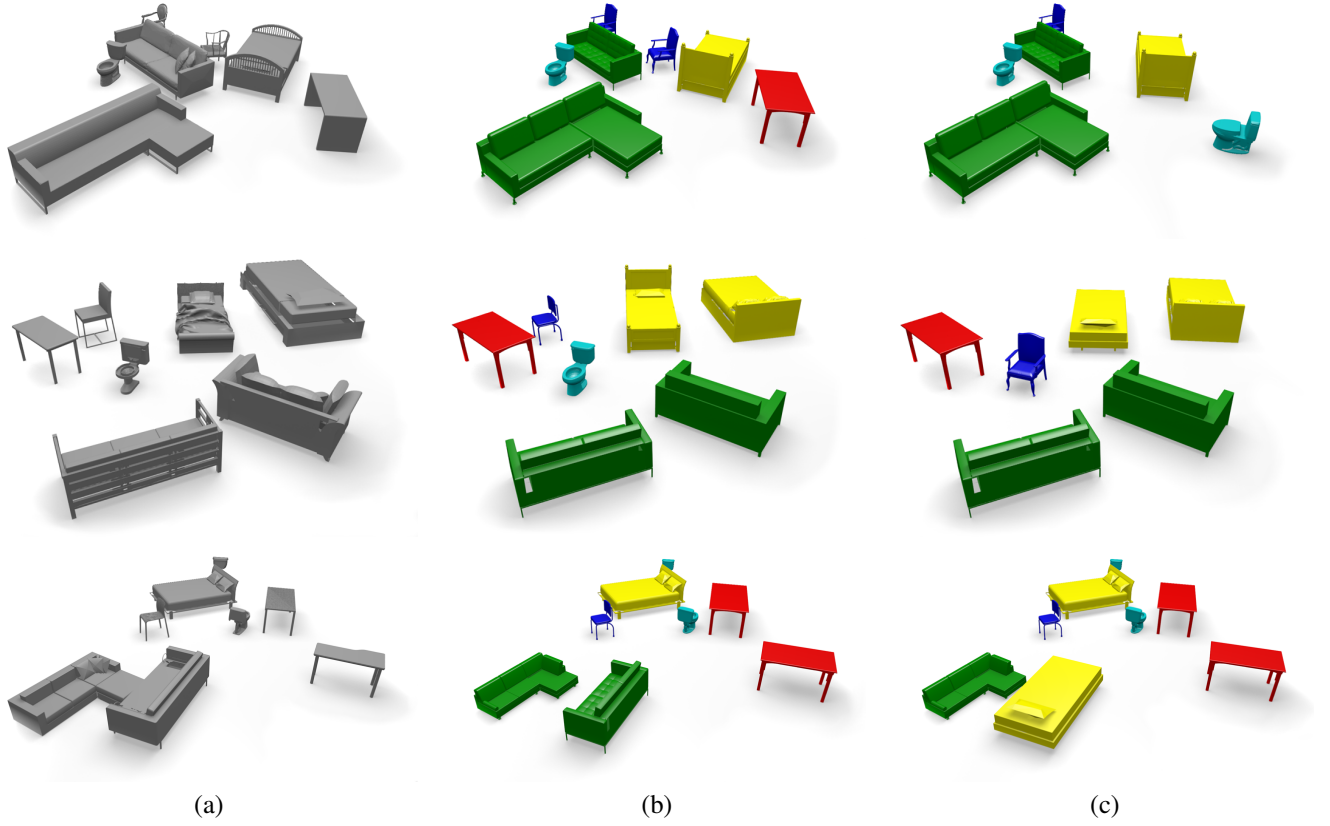


Figure 2: **Results on synthetic scenes.** Column (a) shows the original scene; (b) shows the output of ASIST for $N_{out} = 5$; (c) shows the output of ASIST for $N_{out} = 1$. Each row shows the synthetic scene in gray; and the output of the ASIST algorithm, in which the objects to be inserted are rendered in color-coded scheme according to object class (blue for chairs, red for tables, cyan for toilets, green for sofas, and yellow for beds). The two colored columns (b) and (c) represent the two ASIST configurations with $N_{out} = 5$ and $N_{out} = 1$ respectively.

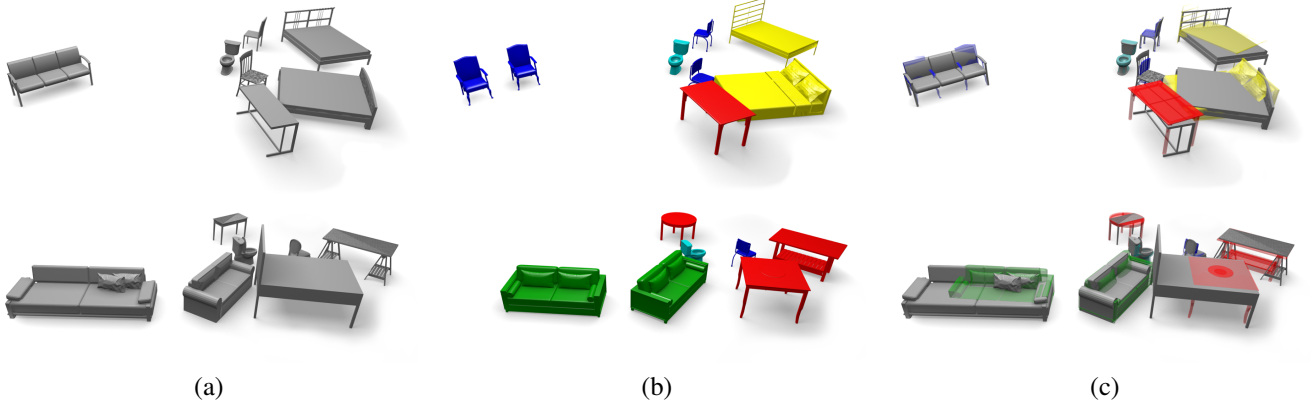


Figure 3: **Failure cases.** Column (a) shows the original scene; (b) shows the output of our algorithm after five iterations; (c) shows the scene and the algorithm's output superimposed.

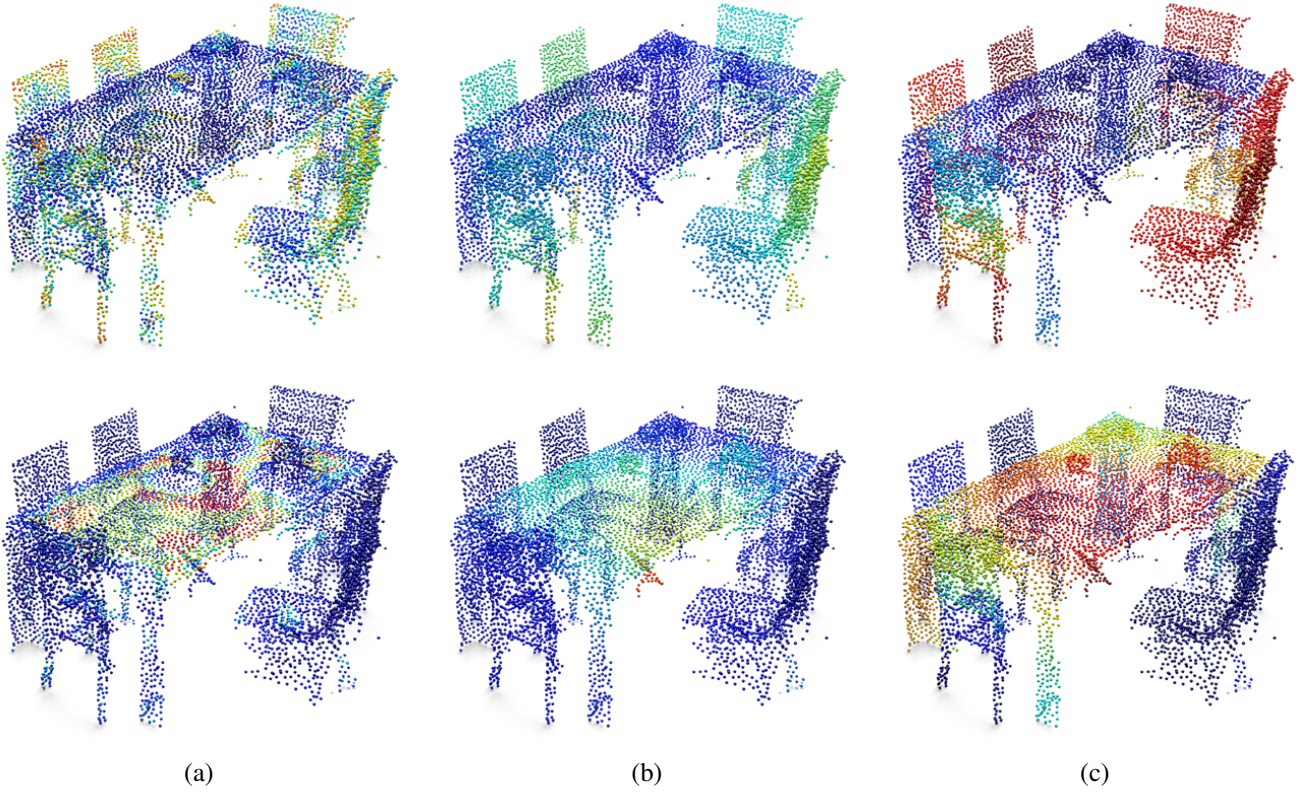


Figure 4: **Pipeline illustration, Part 1.** Per point semantic class segmentation at different steps of ASIST, for the chair class (top row) and table class (bottom row). (a) Initial (forest) labeling; (b) after spectral representation; (c) final segmentation.



Figure 5: **Pipeline illustration, Part 2.** Exemplar locations: (a) at initialization; (b) after a single iteration; (c) after 3 iterations; and (d) after 5 iterations.

Results. In this experiment we present the result of ASIST on six scanned scenes, with floors removed. As can be seen in Figures 4.3-4.3 we get a perfect result in terms of precision and recall on five out of the six scenes. We now elaborate on these results as well as the interesting failure case shown in Figure 4.3.

As can be seen in Figure 4.3, the scans obtained using the Kinect sensor are of a reasonable resolution. Nevertheless they contain a non negligible amount of clutter (see the top of the dining table for example) and noise (see the short coffee table in the second row). These are handled though our clutter-class and the exemplars’ geometry. Recall that due to the alternating fashion of the algorithm the registration, while relying on the segmentation, also aids in improving it.

Taking a closer look at the first example shown in Figure 4.3, the clutter on the table can clearly be seen, while additionally some of the chairs are missing their legs. ASIST handles these issues and returns the correct result as a result of its unified approach. The second example shows how ASIST successfully deals with large occlusions: a large part of the sofa and table are missing from the scan. The third example also contains serious occlusions. It is evident in this case that ASIST returns a rectangular table rather than a round one and chairs that are less reclined than the ones in the scene. This is due to the lack of a better exemplar in the database; nevertheless the result is faithful to the semantics of the scene. The fourth example emphasizes how ASIST handles large and crowded scenes with many different objects and classes, where again, some of the scanned objects have missing parts such as legs and seats. Nevertheless, all objects are recovered correctly with minor geometric errors such as the length of the couch and the leftmost table.

Figure 4.3 shows an example of a scene scanned with the Google Tango sensor. The scan resolution is much lower compared to the Kinect scenes with about half of the spatial resolution. Large parts of the scan are missing and those present are extremely noisy. Yet, these acquisition imperfections are still robustly handled by our algorithm.

In Figure 4.3 we present a failure case of ASIST on a scene scanned using Google Tango. The errors can be divided into two types: geometric, such as the returned table being rectangular instead of round and a few missing chairs at the top right corner of the scene; and semantic, such as the chairs on the left side of the scene that were replaced by sofas. Although this is a semantic error the two class types are relatively close semantically and this sort of error could probably be acceptable in some VR applications.

4.4. Comparison to Li et al. [3]

In this section we compare the performance of ASIST to the algorithm introduced in [3]. Unfortunately, the authors provide no basis for quantitative comparison, due to the lack of annotated scenes and code. However, some of the scans were released along with their reconstructions, and we were therefore able to provide a qualitative comparison. Figure 4.4 shows a comparison between the performance of ASIST (column (b)) and the published reconstruction of Li et al. (column (d)) on two different scenes. For both scenes it can be seen that ASIST

achieves a perfect result in terms of both semantic and geometric recall and precision and thus outperforms the results of Li et al. In particular, it can be seen in the top scene that the solution of Li et al. misses one of the six chairs, while in the bottom scene it misses two of the four chairs. These errors occur at quite challenging parts of the scans: the scanned chairs are missing significant chunks due to occlusions. Nevertheless, ASIST still performs well, inserting chairs appropriately.

Furthermore, it is apparent that the reconstruction stays loyal to the geometry and the pose of the objects in the scene (column (c)). A close look at the results of both algorithms shows that in some of the chairs, for instance, the reconstructed object is a four-legged chair whereas the object in the scene is a swivel chair. In ASIST, this is mainly due to the lack of diversity of exemplars. The set of exemplars we picked contains a single exemplar representing a “swivel chair” with a high backrest, hence for objects such as the second chair from the right and the leftmost one, a lower energy value can be obtained by choosing a chair with higher geometric error in the legs area than in the backrest area.

4.5. Comparison to Nan et al. [2]

A quantitative comparison of ASIST to the algorithm proposed by Nan et al. [2] is reported in Table 3, and a qualitative comparison is shown in Figure 4.5. Since code was neither published nor was made available upon request, we evaluated the performance on the 18 scene dataset that was released by the authors, and compared it against the score derived from their published reconstructions. In order to evaluate precision and recall we manually annotated bounding boxes of chairs, tables and sofas in the scenes. We intend to release this data to the research community.

Observe that while semantically, Nan et al. slightly outperforms ASIST in the majority of cases, the results are comparable both quantitatively and qualitatively (see Table 3). Figure 4.5 shows two of Nan’s published scenes (column (a)), the result of ASIST (column (b)) and Nan’s result (column (d)). In the top row one can observe a semantic failure case by Nan where the algorithm replaces the adjacent chairs with sofas. ASIST, on the other hand, finds the individual chairs correctly. The bottom row shows a case where Nan’s method performs well except for replacing two adjacent chairs with a high table. In that case, ASIST finds all the chairs in the scene correctly, but erroneously replaces the coffee table with a couple of chairs as well.

It is worthwhile noting that while Nan et al. allows deformations of the objects used for the reconstruction, ASIST uses its available exemplars as is. This results in a match which, while performing quite well in terms of precision and recall, deviates geometrically from the ground truth more than Nan’s method. This provides the inspiration for one of the main directions for future research, as discussed in detail in Section 5.

5. Discussion

We have seen that the proposed ASIST algorithm is quite effective in achieving semantically invariant scene transforma-

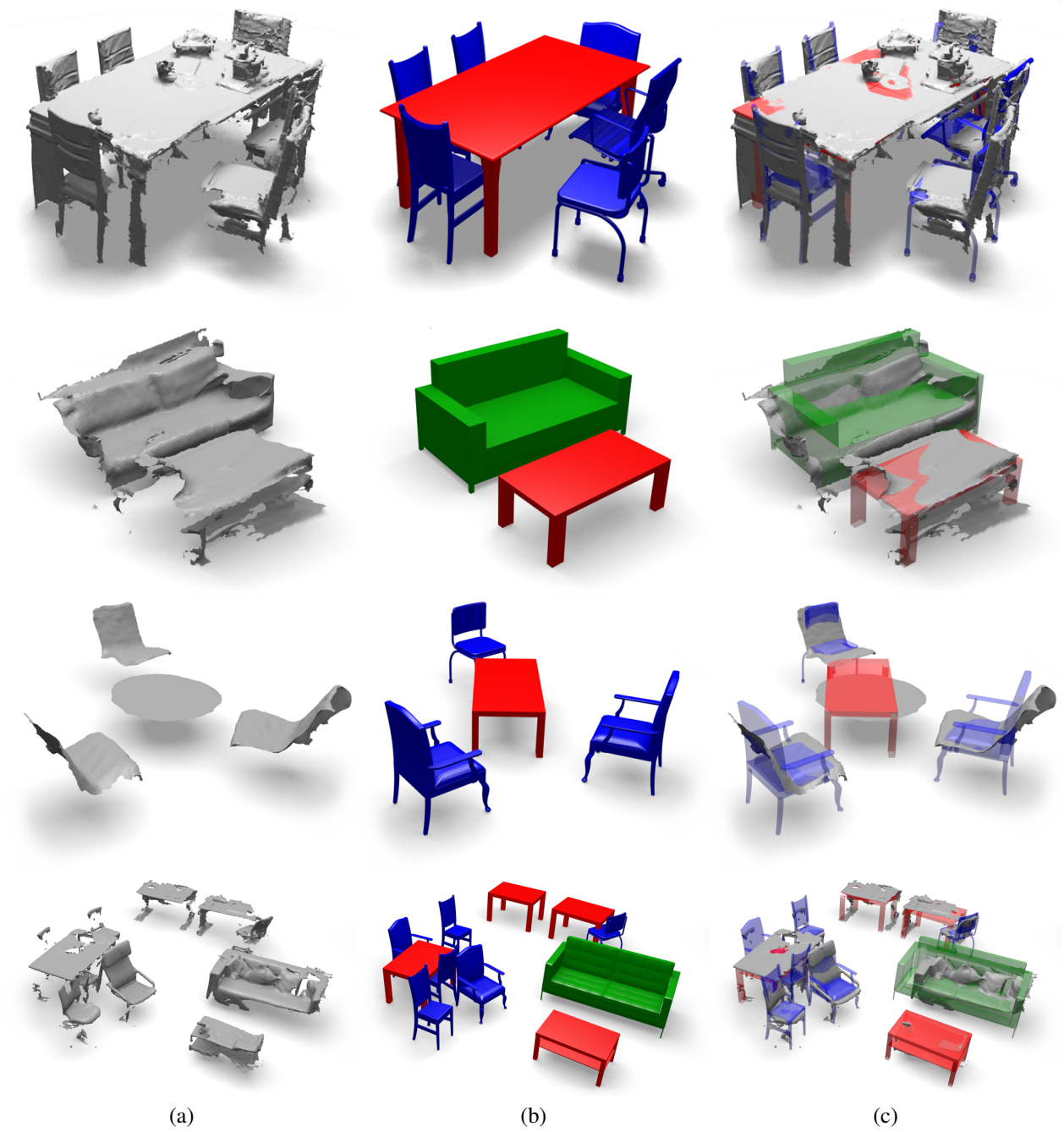


Figure 6: ASIST using a Kinect sensor. Column (a) scanned scene; (b) ASISTs' result; (c) overlay of the previous two.

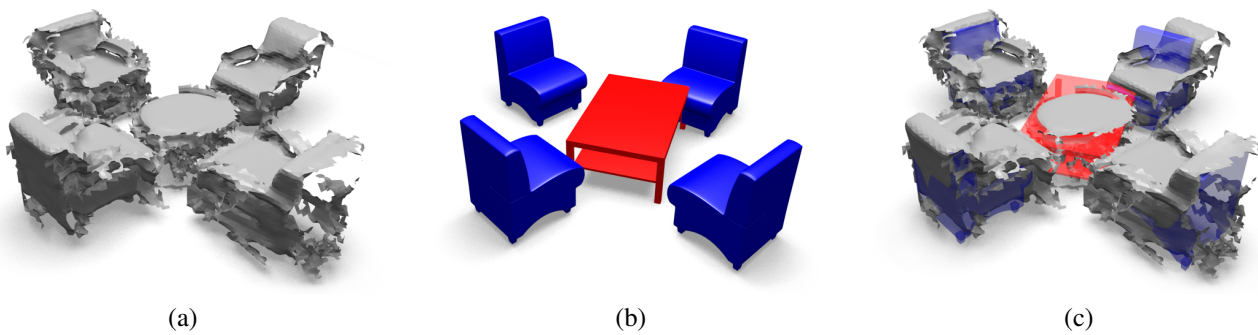


Figure 7: ASIST using a Google Tango sensor. Column (a) scanned scene; (b) ASISTs' result; (c) overlay of the previous two.

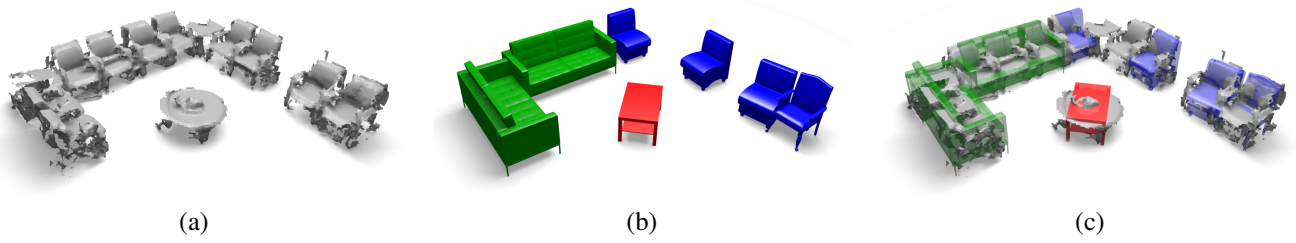


Figure 8: ASISTs' failure on a scene scanned with Google Tango. (a) scanned scene; (b) ASIST result; (c) overlay of the previous two.



Figure 9: Scenes from Li *et al.* [3]. (a) Scanned scene; (b) ASIST result; (c) result overlay; (d) Li *et al.* result, shown in pink. The images of the top example are presented diagonally while the images of the bottom example are presented in a horizontal row.

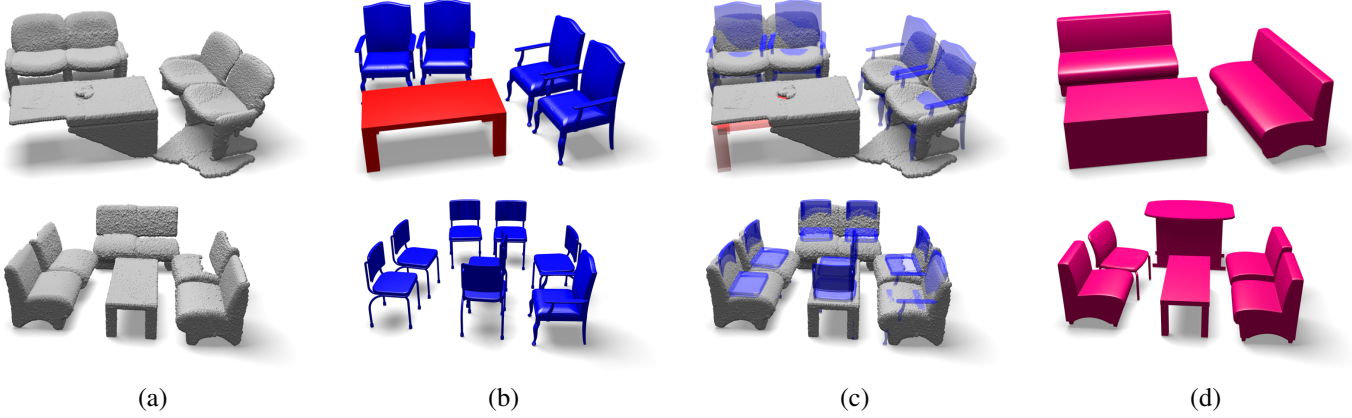


Figure 10: Scenes from Nan *et al.* [2]. (a) Scanned scene; (b) ASIST result; (c) result overlay; (d) Nan *et al.* result, shown in pink.

	chair	table	geo
Precision			
ASIST	0.91	0.8	0.93
Nan <i>et al.</i>	0.97	0.81	0.93
Recall			
ASIST	0.93	0.89	0.97
Nan <i>et al.</i>	0.91	1	0.93
F_1 score			
ASIST	0.92	0.84	0.95
Nan <i>et al.</i>	0.94	0.89	0.93

Table 3: Performance comparison with Nan *et al.* [2]. Semantic and geometric precision, recall and F_1 scores for all scenes published in [2].

tions on a variety of different sources of data. However, in examining the failure cases, it becomes apparent that there are two elements of the algorithm which can lead ASIST astray.

The first element is the cell classification. We have observed that the random forest classifier has reasonable performance, but does not produce extraordinarily accurate results. It turns out that this reasonable level of performance is sufficient for ASIST in many cases, as the other energy terms pull the algorithm towards the correct solution. However, we observed that in many failure cases, it was the forest that provided a signal which was too weak, or even incorrect. The semantic data term was consequently uninformative or incorrect, and the other energy terms could not adequately compensate.

A potential solution to this problem is readily apparent. Much work has been done on object recognition and detection; while the random forest algorithm we rely on is simple to implement, it is naïve and has considerably lower performance than current state-of-the-art recognition algorithms, many of which are based on convolutional neural networks and other deep learning techniques. For example, we might expect that using the detection algorithm introduced in Gupta *et al.* [10], the overall accuracy of the ASIST algorithm would improve. One advantage of the way the ASIST pipeline is built is that plugging in such a state-of-the-art classifier is straightforward.

The second element which negatively affects ASIST’s per-

formance relates to the registration step, i.e., the choice of the per exemplar transformations T_e . In particular, in the current implementation we restrict ourselves to rigid transformations. However, recall that our set of exemplars \mathcal{E} is finite, and relatively small in practice. Thus, it is often the case that an object in the scene will not perfectly match an exemplar, even if the best possible rigid transformation is chosen. Consequently, it is sometimes the case that the resulting match is quite inaccurate.

The issue is that the set of rigid transformations is too restrictive. A remedy is to broaden the set of transformations to include scaling, possibly anisotropic. Further afield, one might consider various classes of non-rigid transformations, for example of the type described in [26]. By expanding the set of transformations, one would expect more accurate matches with the scene. And while the computation due to a broader set of transformations might be more expensive, this could be offset by the need to use fewer exemplars in order to achieve accurate matching.

6. Conclusions

We have presented the ASIST algorithm for computing semantically invariant scene transformations. Due to a unified formulation of semantic segmentation and object replacement based on the optimization of a single objective, ASIST solves both problems simultaneously via an iterative algorithm. The method has been shown to achieve a high level of accuracy on datasets of both synthetic scenes and fused scans, as well as comparable performance to recently published competitor methods [2, 3] on their own data.

References

References

- [1] S. Izadi, D. Kim, O. Hilliges, D. Molyneaux, R. Newcombe, P. Kohli, J. Shotton, S. Hodges, D. Freeman, A. Davison, et al., Kinectfusion: real-time 3d reconstruction and interaction using a moving depth camera, in: Proceedings of the 24th annual ACM symposium on User interface software and technology, ACM, 2011, pp. 559–568.
- [2] L. Nan, K. Xie, A. Sharf, A search-classify approach for cluttered indoor scene understanding, ACM Trans. Graph. 31 (6) (2012) 137:1–137:10.

- [3] Y. Li, A. Dai, L. Guibas, M. Nießner, Database-assisted object retrieval for real-time 3d reconstruction, in: *Computer Graphics Forum*, Vol. 34, 2015.
- [4] S. Gupta, P. Arbeláez, R. Girshick, J. Malik, Aligning 3d models to rgb-d images of cluttered scenes, in: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2015.
- [5] R. Salas-Moreno, R. Newcombe, H. Strasdat, P. Kelly, A. Davison, Slam++: Simultaneous localisation and mapping at the level of objects, in: *Computer Vision and Pattern Recognition (CVPR)*, 2013 IEEE Conference on, 2013, pp. 1352–1359.
- [6] S. Song, J. Xiao, Sliding shapes for 3d object detection in depth images, in: *Computer Vision–ECCV 2014*, Springer, 2014, pp. 634–651.
- [7] Project Tango, <https://www.google.com/atap/project-tango/>.
- [8] B. Drost, M. Ulrich, N. Navab, S. Ilic, Model globally, match locally: Efficient and robust 3d object recognition., in: *CVPR*, IEEE Computer Society, 2010, pp. 998–1005.
- [9] T. Malisiewicz, A. Gupta, A. Efros, et al., Ensemble of exemplar-svms for object detection and beyond, in: *Computer Vision (ICCV)*, 2011 IEEE International Conference on, IEEE, 2011, pp. 89–96.
- [10] S. Gupta, R. Girshick, P. Arbelaez, J. Malik, Learning rich features from RGB-D images for object detection and segmentation, in: *ECCV*, 2014.
- [11] R. Karimi Mahabadi, C. Hane, M. Pollefeys, Segment based 3d object shape priors, in: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 2838–2846.
- [12] C. Hane, C. Zach, A. Cohen, R. Angst, M. Pollefeys, Joint 3d scene reconstruction and class segmentation, in: *Computer Vision and Pattern Recognition (CVPR)*, 2013 IEEE Conference on, IEEE, 2013, pp. 97–104.
- [13] C. Hane, N. Savinov, M. Pollefeys, Class specific 3d object shape priors using surface normals, in: *Computer Vision and Pattern Recognition (CVPR)*, 2014 IEEE Conference on, IEEE, 2014, pp. 652–659.
- [14] D.-Y. Chen, X.-P. Tian, Y.-T. Shen, M. Ouhyoung, On visual similarity based 3d model retrieval, in: *Computer graphics forum*, Vol. 22, Wiley Online Library, 2003, pp. 223–232.
- [15] T. F. et al, Shapenet, <http://shapenet.cs.stanford.edu> (2015).
- [16] Z. Wu, S. Song, A. Khosla, F. Yu, L. Zhang, X. Tang, J. Xiao, 3d shapenets: A deep representation for volumetric shapes, in: *IEEE Conference on Computer Vision and Pattern Recognition, CVPR*, 2015, pp. 1912–1920.
URL <http://dx.doi.org/10.1109/CVPR.2015.7298801>
- [17] S. Song, S. P. Lichtenberg, J. Xiao, Sun rgb-d: A rgb-d scene understanding benchmark suite, 2015.
- [18] A. Criminisi, J. Shotton, E. Konukoglu, Decision forests: A unified framework for classification, regression, density estimation, manifold learning and semi-supervised learning, Now, 2012.
- [19] O. Litany, A. M. Bronstein, M. M. Bronstein, Putting the pieces together: Regularized multi-part shape matching, in: *Computer Vision - ECCV 2012. Workshops and Demonstrations - Florence, Italy, October 7-13, 2012, Proceedings, Part I*, 2012, pp. 1–11.
- [20] S. Rusinkiewicz, M. Levoy, Efficient variants of the icp algorithm, in: *3-D Digital Imaging and Modeling*, 2001. *Proceedings. Third International Conference on*, IEEE, 2001, pp. 145–152.
- [21] R. Chartrand, W. Yin, Iteratively reweighted algorithms for compressive sensing, in: *Acoustics, speech and signal processing*, 2008. *ICASSP 2008. IEEE international conference on*, IEEE, 2008, pp. 3869–3872.
- [22] K. Fukunaga, L. D. Hostetler, The estimation of the gradient of a density function, with applications in pattern recognition, *Information Theory, IEEE Transactions on* 21 (1) (1975) 32–40.
- [23] F. R. Chung, *Spectral graph theory*, Vol. 92, American Mathematical Soc., 1997.
- [24] M. Spagnuolo, M. Bronstein, A. Bronstein, A. Ferreira, et al., Parallelized algorithms for rigid surface alignment on gpu.
- [25] M. Everingham, L. Van Gool, C. K. Williams, J. Winn, A. Zisserman, The pascal visual object classes challenge 2012 (voc2012) results (2012).
- [26] H. Lombaert, L. Grady, X. Pennec, N. Ayache, F. Chérier, Spectral log-demons: diffeomorphic image registration with very large deformations, *International journal of computer vision* 107 (3) (2014) 254–271.